

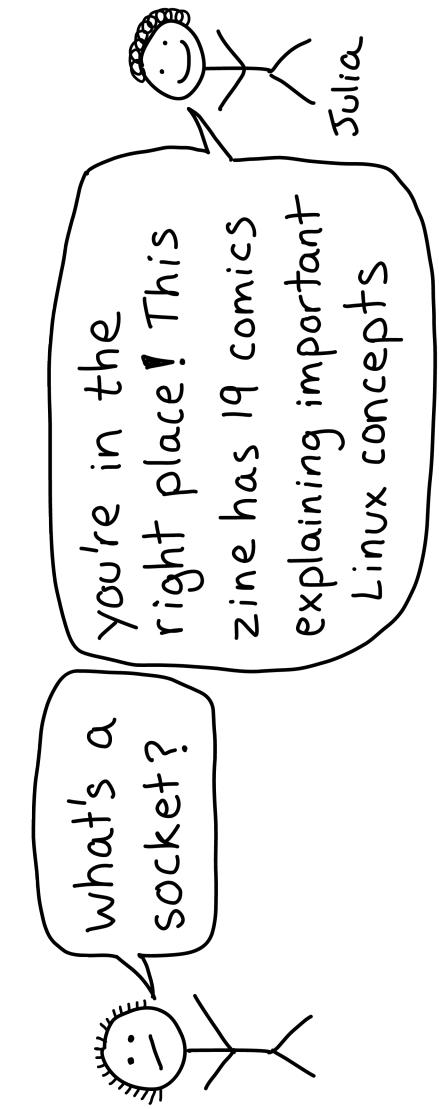
CC-BY-NC-SA computer wizard industries 2018

love this?

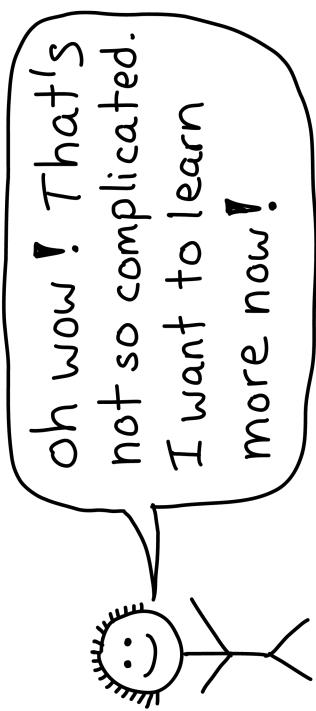
find more awesome zines at

→ jvns.ca/zines ←





... 5 minutes later ...



by Julia Evans
<https://jvns.ca>
twitter.com/b0rk

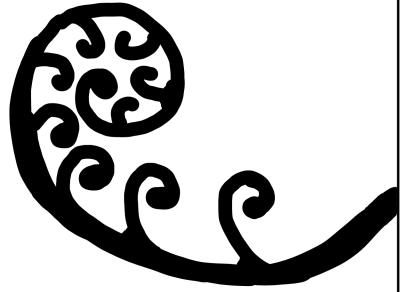
Want to learn more?
I highly recommend
this book →

THE LINUX PROGRAMMING INTERFACE

MICHAEL KERRISK

Every chapter is a readable,
short (usually 10-20 page)
explanation of a Linux system.
I used it as a reference
constantly when writing
this zine.

I ♡ it because even though
it's huge and comprehensive
(1500 pages!), the chapters
are short and self-contained
and it's very easy to pick it
up and learn something.



man page sections 22

man pages are split up
into 8 sections

① ② ③ ④ ⑤ ⑥ ⑦ ⑧

\$ man 2 read

means "get me the man page
for read from section 2"

There's both

→ a program called "read"
→ and a system call called "read"

\$ man 1 read

gives you a different man page from

\$ man 2 read

If you don't specify a section, man will
look through all the sections & show
the first one it finds

man page sections

① programs ② system calls

\$ man grep

\$ man sendfile

\$ man ls

\$ man ptrace

③ C functions ④ devices

\$ man printf

\$ man null

\$ man fopen

for /dev/null docs

⑤ file formats ⑥ games

\$ man sudoers

for /etc/sudoers

\$ man proc

files in /proc!

\$ man sl

not super useful.

⑦ miscellaneous ⑧ sysadmin programs

\$ man 7 pipe

explains concepts!

\$ man 7 symlink

is funny if you have

\$ man apt

\$ man chroot

♥ Table of contents ♥

unix permissions	4	sockets	10	virtual memory	17
/proc	5	unix domain sockets	11	shared libraries	18
system calls	6	processes	12	copy on write	19
signals	7	threads	13	page faults	20
file descriptors	8	floating point	14	mmap	21
pipes	9	file buffering	15	man page sections	22
memory allocation	16				

unix permissions

4

There are 3 things you can do to a file

→ read write execute

$r w -$ $r w -$ $r --$
↑ ↑ ↑
bork (user) staff (group) ANYONE
can read & write can read & write can read

`ls -l file.txt` shows you permissions
Here's how to interpret the output:

File permissions are 12 bits

setuid setgid ↓ user group all = 110 100 100
000 110 110 100
sticky rwx rwx rwx
For files:
 r = can read
 w = can write
 x = can execute
For directories it's approximately:
 r = can list files
 w = can create files
 x = can cd into & access files

File permissions are 12 bits

110 in binary is 6
 $So \ 110 \ 100 \ 100$
 $= \ 6 \ 4 \ 4$
chmod 644 file.txt + means change the permissions to:
 $r w - \ r -- \ r --$
Simple!

Setuid affects executables

\$ ls -l /bin/ping
rws r-x r-x root root
this means ping always runs as root

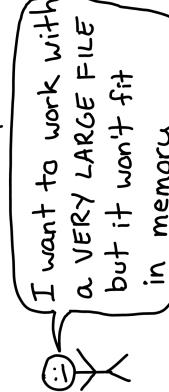
setgid does 3 different unrelated things for executables, directories, and regular files

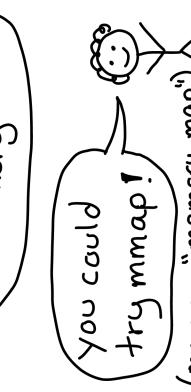
 unix! it's a long story why???

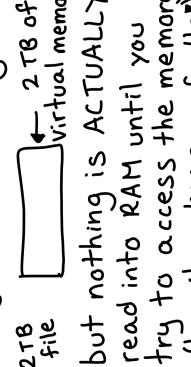
mmap

21

What's mmap for?

 I want to work with a VERY LARGE FILE but it won't fit in memory

 You could try mmap!
(mmap = "memory map")

load files lazily with mmap
When you mmap a file, it gets mapped into your program's memory

but nothing is ACTUALLY read into RAM until you try to access the memory (how it works: page faults!)

how to mmap in Python

import mmap

f = open("HUGE.txt", "r")

mm = mmap.mmap(f.fileno(), 0)

this won't read the file from disk!

Finishes ~instantly.

print(mm[-1000:])

this will read only the last 1000 bytes!

anonymous memory maps

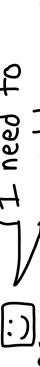
→ not from a file (memory set to 0 by default)

with MAP_SHARED, you can use them to share memory with a subprocess!

dynamic linking uses mmap

 I need to use libc.so.6

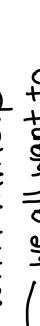
Standard library

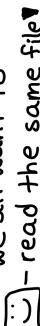
 you too eh? no problem

I always mmap, so that file is probably loaded into memory already

sharing big files with mmap

 we all want to read the same file!

 no problem!

 mmap

Even if 10 processes mmap a file, it will only be read into memory once.

Page faults

every Linux process has a page table

* page table *

virtual memory address

virtual memory address	physical memory address
0x19723000	0x1422000
0x19724000	0x1423000
0x1524000	not in memory
0x1844000	0x1a000 read only

"not in memory" usually means the data is on disk!

virtual memory



Having some virtual memory that is actually on disk is how swap and mmap work

some pages are marked as either

* read only

* not resident in memory

when you try to access a page that's marked "not in memory", that triggers a ! page fault!

- your program stops running
- the MMU sends an interrupt
- Linux kernel code to handle the page fault runs
- I'll fix the problem and let your program keep running

how swap works

- ① run out of RAM
 - ② Linux saves some RAM data to disk
 - ③ mark those pages as "not resident in memory" in the page table, not resident
 - ④ When a program tries to access the memory there's a ! page fault!
 - ⑤ Linux → time to move some data back to RAM!
 - ⑥ if this happens a lot your program gets VERY SLOW
- I'm always waiting for data to be moved in & out of RAM

an amazing directory: /proc 5

Every process on Linux has a PID (process ID) like 42.

In /proc/42, there's a lot of **VERY USEFUL** information about process 42

/proc/PID/fd

Directory with every file the process has open!

Run \$ ls -l /proc/42/fd to see the list of files for process 42.

These symlinks are also magic & you can use them to recover deleted files ♥

- The kernel's current stack for the process. Useful if it's stuck in a system call
- List of process's memory maps. Shared libraries, heap, anonymous maps, etc.

/proc/PID/stack

Look at

man proc

for more information!

/proc/PID/exe symlink to the process's binary magic: works even if the binary has been deleted!

/proc/PID/status Is the program running or asleep? How much memory is it using? And much more!

And :more:

Look at

man proc

for more information!

System calls

The Linux kernel has code to do a lot of things

- read from a hard drive
- make network connections
- create new process
- kill a process
- change file permissions
- Keyboard drivers

every program uses system calls

- I use the 'open' syscall to open files
- me too!
- me three!
- C program
- Python program

your program doesn't know how to do those things

- TCP? dude I have no idea how that works
- NO I do not know how the ext4 filesystem is implemented I just want to read some files

and every system call has a number (eg chmod is #90 on x86-64)

So what's actually going on when you change a file's permissions is

- run syscall #90 program with these arguments

please write to this file
(switch to running kernel code)

done! I wrote 1097 bytes! Linux
<program resumes>

please write to this file
(switch to running kernel code)

copy on write

On Linux, you start new processes using the fork() or clone() system call

calling fork gives you a child process that's a copy of you

Parent child

so Linux lets them share physical RAM and only copies the memory when one of them tries to write.

I'd like to change that memory

ok I'll make you your own copy! Linux

copying all that memory every time we fork would be slow and a waste of RAM

often processes call exec right after fork which means they don't use the parent process's memory basically at all!

when a process tries to write to a shared memory address

- ① there's a ≈ page fault =
- ② Linux makes a copy of the page & updates the page table
- ③ the process continues, blissfully ignorant

It's just like I have my own copy!

same ← RAM

Shared libraries

18

Most programs on Linux use a bunch of C libraries some popular libraries:

`openssl`

`(for SSL!)`

`sqlite`

`(embedded db!)`

`libpcre`

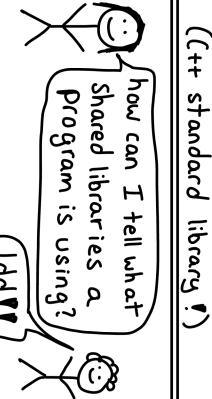
`(regular expressions!)`

`zlib`

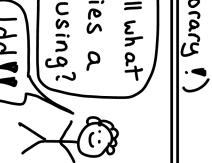
`(gzip!)`

`libstdc++`

`(C++ standard library!)`



how can I tell what shared libraries a program is using?



\$ ldd /usr/bin/curl
libz.so.1 => /lib/x86_64...
libresolv.so.2 => ...
libc.so.6 => ...
+ 34 more !!

There are 2 ways to use any library

① Link it into your binary

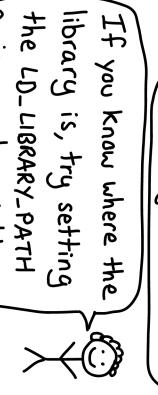
`[your code] z! lib sqlite`

big binary with lots of things!

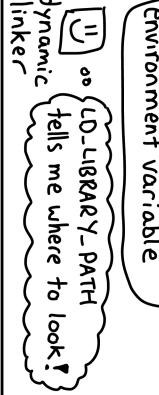
② Use separate shared libraries

`[your code] [z! lib] [sqlite]`

← all different files



If you know where the library is, try setting the LD_LIBRARY_PATH environment variable



`[!] oo [LD-LIBRARY-PATH dynamic linker]`
`tells me where to look!`

Programs like this:

`[your code] z! lib sqlite`

are called "statically linked"

and programs like this:

`[your code] [z! lib] [sqlite]`

are called "dynamically linked"

Where the dynamic linker looks

① DT_RPATH in executable

③ DT_RUNPATH in executable

④ /etc/ld.so.cache

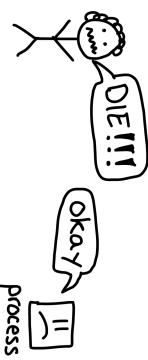
(run ldconfig -p to see contents)

⑤ /lib, /usr/lib

Signals

7

If you've ever used  kill you've used signals



You can send signals yourself with the kill system call or command

SIGINT ctrl-C } various levels of
SIGTERM kill } "die"
SIGKILL kill -q }

SIGHUP kill -HUP often interpreted as "reload config", eg by nginx

Every signal has a default action, one of:

ignore

kill process

kill process AND make core dump file

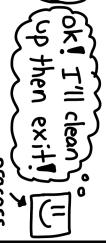
stop process

resume process

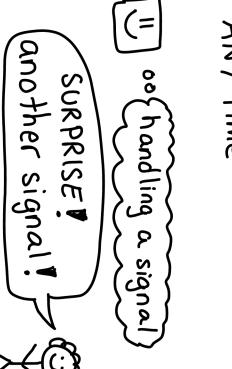


SIGTERM (terminate) ok! I'll clean up then exit!

process

exceptions:
SIGSTOP & SIGKILL can't be ignored
got → 

Signals can be hard to handle correctly since they can happen at ANY time



`[!] oo [handling a signal]`

SURPRISE!
another signal!

file descriptors

8

Unix systems use integers to track open files

`Open foo.txt`

Okay! that's file #7 for you.

these integers are called file descriptors

`lsof (list open files) will show you a process's open files`

`$ ls -p 4242 ← PID were interested in FD NAME /dev/pts/ty1`

`1 /dev/pts/ty1`

`2 pipe:29174`

`3 /home/bark/awesome.txt`

`5 /tmp/`

`FD is for file descriptor`

When you read or write to a file/pipe/network connection you do that using a file descriptor

`connect to google.com`

`ok! fd is 5!`

`OS`

`write GET / HTTP/1.1 to fd #5 done!`

file descriptors can refer to:

- files on disk
- pipes
- sockets (network connections)
- terminals (like xterm)
- devices (your speaker! /dev/null!)
- LOTS MORE! (eventfd, inotify, signals, epoll, etc etc)

not EVERYTHING on Unix is a file, but lots of things are

virtual memory

17

every program has its own virtual address space

`0x129520 → "puppies"`

`program 1`

`0x129520 → "bananas"`

`program 2`

every time you switch which process is running, Linux needs to switch the page table

`here's the address of process 2950's page table`

`Linux`

`thanks I'll use that now!`

physical memory has addresses 0 - 8GB

but when your program references an address like 0x5c69a2a2

`↑`

that's not a physical memory address!

It's a virtual address

when your program accesses a virtual address

`I'm accessing 0x21000`

`CPU`

`of I'll look that up in the page table and then access the right physical address`

`MMU`

`memory management unit`

`hardware`

Linux keeps a mapping from virtual memory pages to physical memory pages called the "page table"

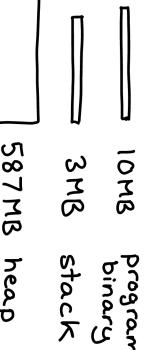
`a "page" is a 4KB or sometimes bigger chunk of memory`

PID	virtual addr	physical addr
1971	0x20000	0x102000
2310	0x20000	0x228000
2310	0x21000	0x9788000

memory allocation

16

Your program has memory



the heap is what your allocator manages

Your memory allocator's interface

malloc (size_t size)
allocate size bytes of memory & return a pointer to it
free (void* pointer)
mark the memory as unused (and maybe give back to the OS)
realloc(void* pointer, size_t size)
ask for more / less memory for pointer
calloc (size_t members, size_t size)
allocate array + initialize to 0

Your memory allocator (malloc) is responsible for 2 things.

THING 1: keep track of what memory is used / free



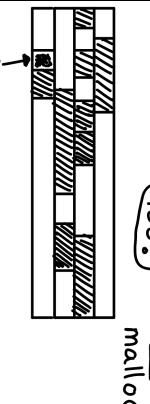
malloc tries to fill in unused space when you ask for memory

your code

can I have 512 bytes of memory?

YES!

malloc



malloc isn't magical! it's just a function!

you can always:

- use a different malloc library like jemalloc or tcmalloc (easy!)
- implement your own malloc (harder)

THING 2: Ask the OS for more memory!

oo oh no I'm being asked for 40 MB and I don't have it

malloc
here you go! OS
can I have 60MB more?

pipes

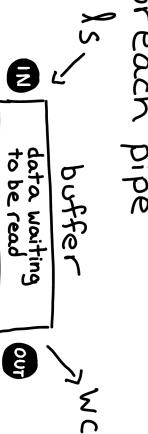
9

Sometimes you want to send the output of one process to the input of another

```
$ ls | wc -l
```

53
53 files!

Linux creates a buffer for each pipe

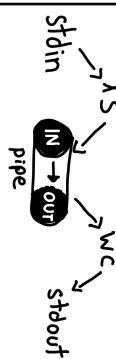


If data gets written to the pipe faster than it's read, the buffer will fill up. When the buffer is full, writes to IN will block (wait) until the reader reads. This is normal & ok!

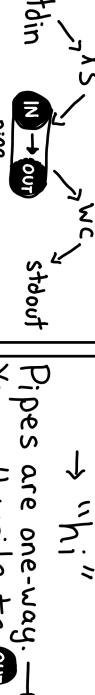
a pipe is a pair of 2 magical file descriptors

① pipe input
IN → OUT

② pipe output
OUT ← IN



what if your target process dies?



If wc dies, the pipe will close and ls will be sent SIGPIPE. By default SIGPIPE terminates your process.

When ls does write(IN, "hi")

wc can read it!
read(OUT)

→ "hi"

Pipes are one-way. → You can't write to OUT.

named pipes

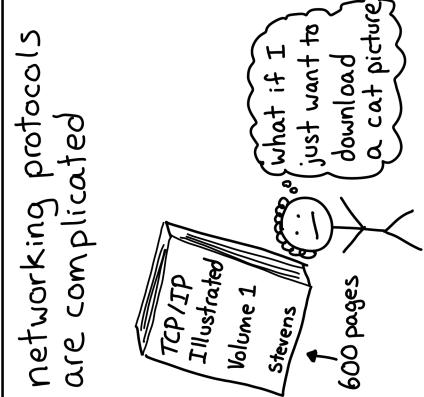
\$ mkfifo my-pipe

This lets 2 unrelated processes communicate through a pipe!

```
f=open("./my-pipe")
f.write("hi\n")
f.readline() ← "hi!"
```

Sockets

10



Networking protocols are complicated
an API called the "socket API" that makes it easier to make network connections (Windows too! :))

you don't need to know how TCP works, Unix I'll take care of it!

Every HTTP library uses sockets under the hood

\$ curl awesome.com
Python: requests.get("yay.us")

oh, cool, I could write a HTTP library too if I wanted. * Neat! * SO MANY edge cases though! :)

here's what getting a cat picture with the socket API looks like:

- ① Create a socket
- fd = socket(AF_INET, SOCK_STREAM)
- ② Connect to an IP/port
- connect(fd, (struct sockaddr_in){ .sin_addr.s_addr = 12.13.14.15, .sin_port = 80 })
- ③ Make a request
- write(fd, "GET /cat.png HTTP/1.1")
- ④ Read the response
- cat-picture = read(fd, ...)

3 Kinds of internet (AF_INET) sockets:

- SOCK_STREAM = TCP
curl uses this
- SOCK_DGRAM = UDP
dig (DNS) uses this
- SOCK_RAW = just let me send IP packets
ping uses this
I will implement my own protocol

file buffering

15

? ? ? I printed some text but it didn't appear on the screen. why?? time to learn about flushing!

On Linux you write to files & terminals with a system call called write

please write "I ❤ cats" to file #1 (stdout)

Okay! Linux

I/O libraries don't always call write when you print

```
printf("I ❤ cats");
```

:) I'll wait for a newline printf before actually writing

This is called buffering and it helps save on syscalls

3 kinds of buffering (defaults vary by library)

- ① None. This is the default for stderr (write after newline). The default for terminals.
- ② Line buffering.
- ③ "full" buffering. (write in big chunks). The default for files and pipes.

flushing

To force your IO library to write everything it has in its buffer right now, call flush!

:) I'll call write right away!

:) no seriously, actually write to that pipe please

floating point

14

a double is 64 bits.

sign exponent fraction
 \uparrow
 \uparrow
 100101 100101 100101 100101 100101 100101 100101 100101

$\pm 2^{\text{E}-1023} \times 1.\text{frac}$

That means there are 2^{64} doubles.
 The biggest one is about 2^{1023} .

doubles get farther apart as they get bigger between 2^n and 2^{n+1} there are always 2^{52} doubles, evenly spaced that means the next double after 2^{60} is $2^{60} + 64 = \frac{2^{60}}{2^{52}}$

$$2^{52} + 0.2 = 2^{52}$$

$1 + \frac{1}{2^{54}} = 1$ ← (the next number after 2^{52} is $2^{52} + 1$)

$$1 + \frac{1}{2^{52}} = 1 \quad \leftarrow \text{infinity is a double}$$

$2^{2000} = \text{infinity}$ ← infinity - infinity ← nan = "not a number"

JavaScript only has doubles (no integers!)

$$> 2^{**53} \\ 9007199254740992 \\ > 2^{**53} + 1 \\ 9007199254740992 \\ \text{same number! uh oh!}$$

doubles are scary and their arithmetic is weird

they're very logical! just understand how they work and don't use integers over 2^{53} in JavaScript

Unix domain sockets !!

unix domain sockets are files.

```
$ file mysock.sock
socket
```

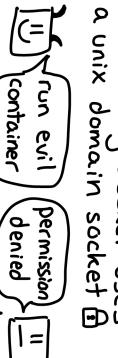
the file's permissions determine who can send data to the socket

advantage 1

lets you use file permissions to restrict access to HTTP/ database services!

chmod 600 secret.sock

This is why Docker uses a unix domain socket



advantage 2

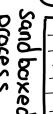
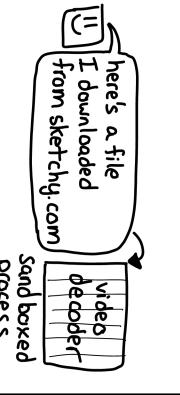
UDP sockets aren't always reliable (even on the same computer). Unix domain datagram sockets are reliable!

And won't reorder!



advantage 3

You can send a file descriptor over a unix domain socket. Useful when handling untrusted input files!



what's in a process?

PID	USER and GROUP who are you running as? julia!	ENVIRONMENT VARIABLES like PATH! you can set them with: \$ env A=val ./program	SIGNAL HANDLERS I ignore SIGHUP! I shut down safely!
WORKING DIRECTORY	Relative paths (./blah) are relative to the working directory! chdir changes it.	PARENT PID PID 1 (init) ↓ PID 147 ↑ PID 129	COMMAND LINE ARGUMENTS see them in /proc/PID/cmdline
MEMORY	heap! stack! shared libraries! the program's binary! mmaped files!	CAPABILITIES I have CAP_PTRACE well I have CAP_SYS_ADMIN	NAMESPACES I'm in the host network namespace I have my own namespace! container process

threads

Threads let a process do many different things at the same time	and they share code calculate-pi find-big-prime-number	why use threads instead of starting a new process? → a thread takes less time to create
process:	I'll write some digits of π to 0x129420 in memory oh that's where I was putting my prime numbers	sharing data between threads is very easy. But it's also easier to make mistakes with threads
	thread 1 I'm calculating ten million digits of π! so fun!	You weren't supposed to CHANGE that data!
	thread 2 I'm finding a REALLY BIG prime number!	thread 1

RESULT: 24 ↪ WRONG.
Should be 25!